# Hybrid enrichment – Anchored phylogenomics

# Collection of High-Throughput Data

- Targeted Amplicon Sequencing or Parallel Tagged Sequencing
- Multiplex PCR
- Massively Parallel Uniplex PCR
- Reduced Representation Library or RAD Sequencing
- Transcriptome Sequencing
- **Hybrid Enrichment**

- **BUT: more data = better data ??**

**Table 1**

Methods of sample preparation for using NGS in phylogeography and phylogenetics.

| Method | Other names or variants | Literature method | Literature examples | Benefits | Drawbacks | Best application |
|---|---|---|---|---|---|---|
| Amplicon sequencing | Multiplex PCR, parallel tagged sequencing | Binladen et al. (2007), Meyer et al. (2008), Tewhey et al. (2009b) | Chan et al. (2010), Griffin et al. (2011), Gunnarsdóttir et al. (2011); Morin et al. (2010), Parks et al. (2009) | Highly targeted. Results in nearly complete data matrices. Needed coverage easy to calculate. Circumvents individual sequencing reactions and phasing nuclear loci compared to Sanger sequencing | Requires PCR of each individual at each locus | Small- to medium-scale projects targeting a limited number of genes |
| Restriction-digest | Double-digest genome reduction, RAD sequencing (RAD-seq), complexity reduction of multilocus sequences (CRoPS), Genotyping by Sequencing (GBS) | Baird et al. (2008), Davey et al. (2011) | Andolfatto et al. (2011), Amaral et al. (2009), Bers et al. (2010), Emerson et al. (2010), Gompert et al. (2010), Hohenlohe et al. (2011), Hyten et al. (2010a,b), Kerstens et al. (2009), Ramos et al. (2009); Sánchez et al. (2009); Van Orsouw et al. (2007); Van Tassell et al. (2008), Wiedmann et al. (2008); Williams et al. (2010) | Broad, random genomic sampling of thousands of independent genomic regions. Requires no prior genomic resources whatsoever | Reproducibility and throughput may be limited by gel extraction step. Not targeted, thus coverage can be difficult to estimate. Null alleles could skew diversity estimates | Intraspecific studies of recent divergence |
| Target enrichment | Sequence capture, targeted resequencing, primer extension capture (PEC) | Albert et al. (2007), Gnirke et al. (2009), Hodges et al. (2007), Okou et al. (2007), Tewhey et al. (2009a), Maricic et al. (2010) | Briggs et al. (2009; Faircloth et al. in press) | Rapid collection of thousands of loci without individual PCR | Requires some prior genomic knowledge, but not necessarily a sequenced genome | Phylogenomics at taxonomic levels above at and above the species level |
| Transcriptome | RNA-seq | Morin et al. (2008); Marioni et al. (2008) | Barbazuk and Schnable (2011), Cánovas et al. (2010), Chepelev et al. (2009); Geraldes et al. (2011), Hittinger et al. (2010) | Can leverage data from expression studies | Skewed read distributions can outstrip coverage, making it difficult to find orthologous loci | Leveraging existing cDNA libraries |

# Comparison of Methods

# Ultraconserved Elements Anchor Thousands of Genetic Markers Spanning Multiple Evolutionary Timescales

BRANT C. FAIRCLOTH[1,*], JOHN E. MCCORMACK[2], NICHOLAS G. CRAWFORD[3], MICHAEL G. HARVEY[2,4], ROBB T. BRUMFIELD[2,4], AND TRAVIS C. GLENN[5]
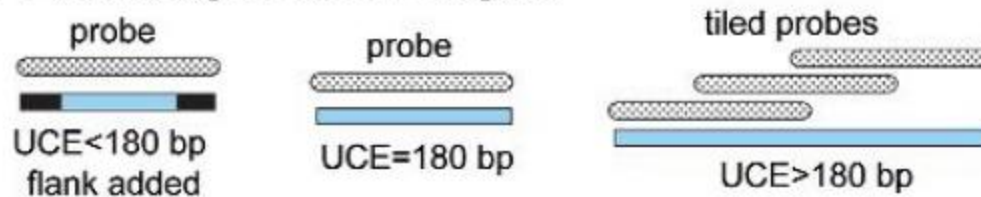
# Anchored Hybrid Enrichment for Massively High-Throughput Phylogenomics

ALAN R. LEMMON[1,*], SANDRA A. EMME[2], AND EMILY MORIARTY LEMMON[2]

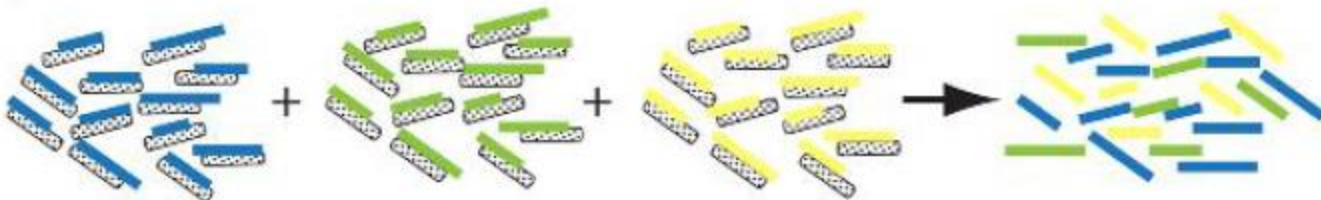# UCE workflow

a) UCEs identified in alignments of birds and lizard

b) Probes designed from UCE regions

c) RNA probes mixed with sheared genomic DNA from non-model organisms

d) Target DNA isolated, enriched, tagged, and pooled for NGS

e) Contigs assembled from NGS reads, aligned to probe, and consensus called for locus

locus *j*

consensus

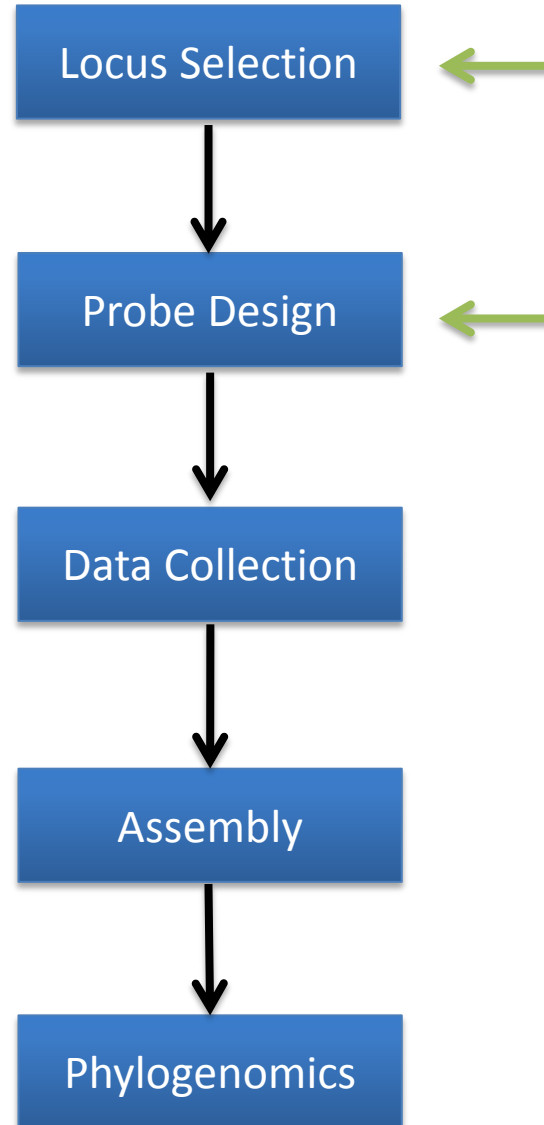f) Consensus loci aligned among species and gene trees estimated for all loci $j_{1 \to n}$

locus *j*          locus *k*          locus *m*

g) Species tree estimated from gene trees

- Locus Selection
  - ~500 "single copy" loci (typically long exons)
  - Conserved element (~20% divergence required)
  - Adjacent to less conserved regions
  - Loci are selected based on broad taxonomic group (e.g., vertebrates)

- Probe Design
  - Incorporate sufficient number of lineages
  - Tile probes across conserved region
  - Goal is to capture ~1500bp regions
  - Probe sets are designed for project-specific clade

# Probe Design Workflow

Probe Design Diagram

# Phases of Hybrid Enrichment for Phylogenomics

Bioinformatics of locus selection — Choose target loci using genomes

↓

Bioinformatics of probe design — Tile probes across target regions

↓

Wet-lab sample processing — Genomic DNA→ raw sequence reads

↓

Bioinformatics of raw data analysis — Raw sequence reads→
phased alleles, alignments, models of evolution, concordance analyses, and phylogenies

# Solution-Phase Hybridization



**GENOMIC SAMPLE**
(Set of chromosomes)

NGS Kit

**GENOMIC SAMPLE (PREPPED)** + **SureSelect HYB BUFFER** + **SureSelect BIOTINYLATED RNA LIBRARY "BAITS"**

**SureSelect Target Enrichment System Capture Process**

= Probes

nature biotechnology

Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted sequencing

Andreas Gnirke[1], Alexandre Melnikov[1], Jared Maguire[1], Peter Rogov[1], Emily M LeProust[2], William Brockman[1,5], Timothy Fennell[1], Georgia Giannoukos[1], Sheila Fisher[1], Carsten Russ[1], Stacey Gabriel[1], David B Jaffe[1], Eric S Lander[1,3,4] & Chad Nusbaum[1]

Hybridization

STREPTAVIDIN COATED MAGNETIC BEADS

Bead capture

UNBOUND FRACTION DISCARDED

Wash Beads and Digest RNA

Amplify

Sequencing

Gnirke et al. 2009

# Hybrid Enrichment

- Disadvantages:
  - Large equipment investment for small operations
    - Equipment required (Bioanalyzer, Covaris or other sonicator, qPCR machine, etc.)
  - Bioinformatic training required for locus selection, probe design, analysis of raw data
  - Substantial investment in reagents
    - Indexes for library preparation
    - Other library preparation reagents
    - Hybrid enrichment kit

# Hybrid Enrichment

- Advantages:
  - Large quantity of data
  - Complete data matrices! Can trim out all missing data and still have more than enough
  - Fast data collection (DNA to phylogeny 2 weeks)
  - Works on degraded samples and ancient DNA
  - DNA starting material (RNA not needed)
  - Can use single probe design (kit) for broad taxonomic group (e.g., Vertebrata)
  - No problem for non-model systems

**a)** Number of Extant Vertebrate Lineages vs. Age (Ma)

Fish
Amphibians
Reptiles
Birds
Mammals

**b)** Vertebrate Lineages Requiring Representation vs. Desired Number of Loci

**c)** Fish, Amphibians, Reptiles, Birds, Mammals

# Own experiences: samples → trees



Requirement: > 1.5- 2.0 micrograms of RNA-free DNA, quantified by Qubit

# Output summary

# Description of the assembled data (example includes nine samples)

nGenes=434
Number of taxa = 9
Number of sites = 711063
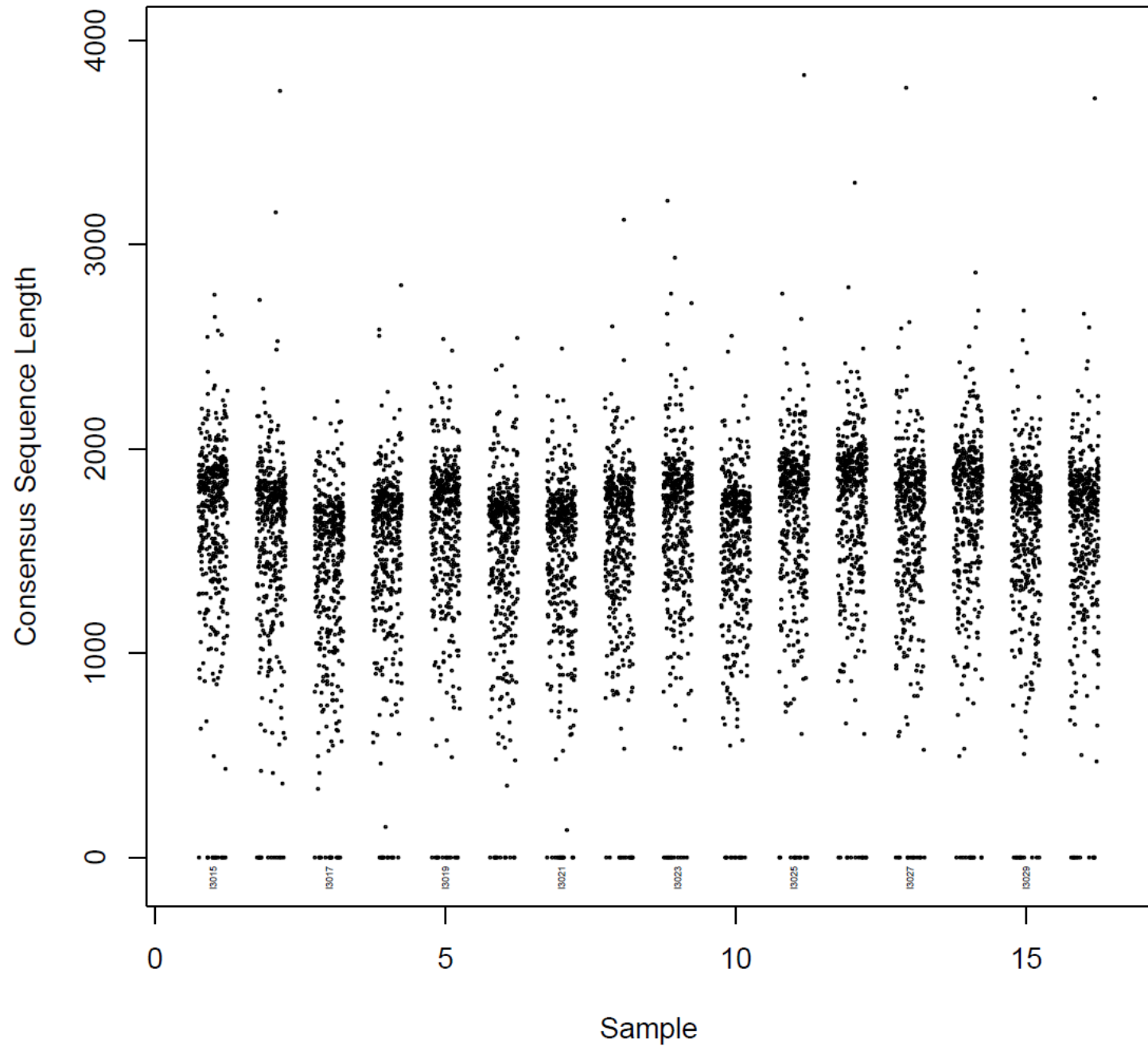Number of variable sites = 16986
Number of informative sites = 4373
Number of characters (total) = 6399567
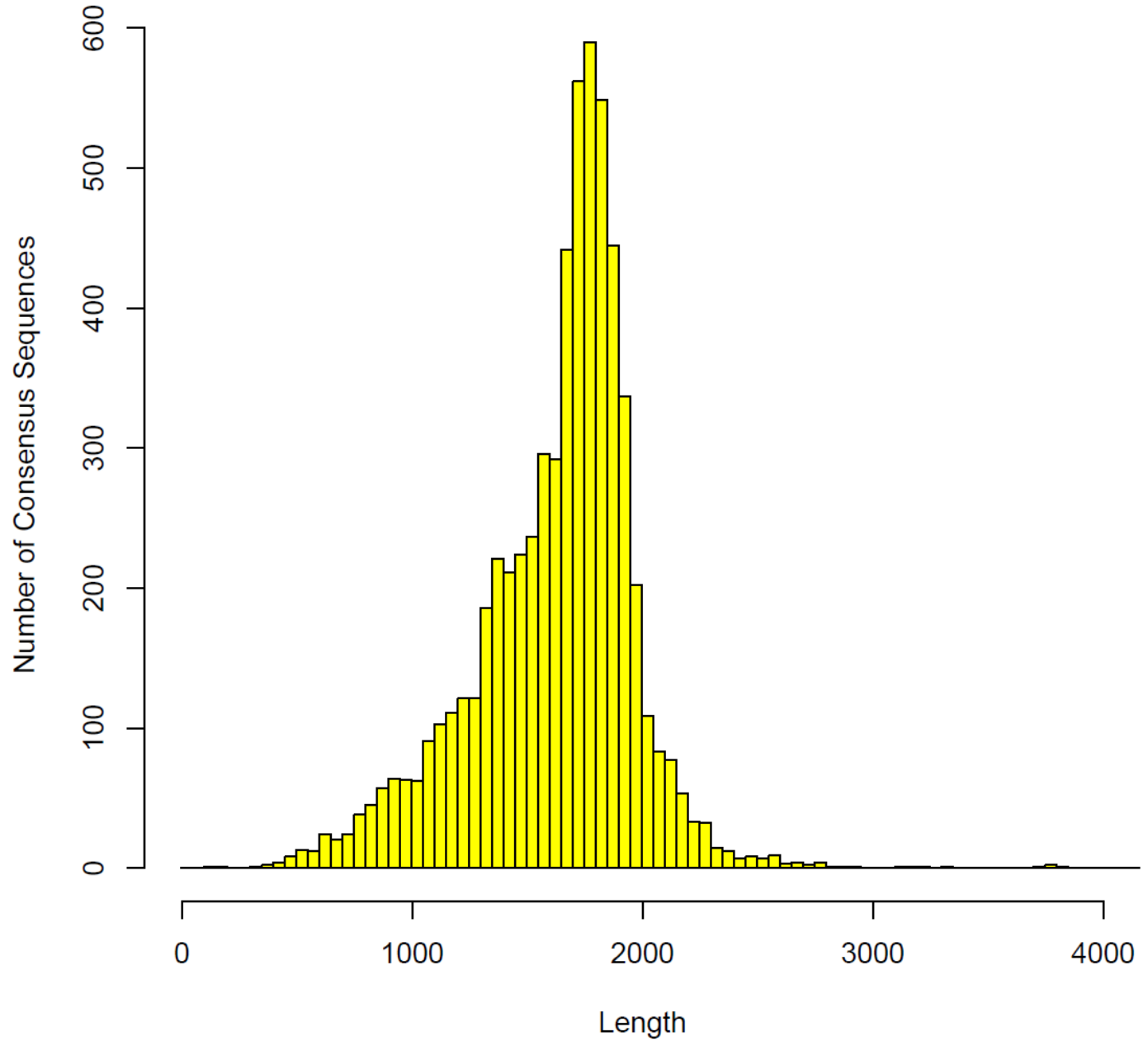% Missing characters (N's and -'s considered only) =1.4503637511725402
% Missing characters (all considered) =1.5992800762926618

Assembly Results − P0041

**Assembly Results − P0041**

**Assembly Results − P0041**

# Approaches to analyze phylogenomic data

**Important issues to consider:**

<span style="color:red">**!!!**</span> amount of missing/ambiguous data

<span style="color:red">**!!!**</span> alignment

<span style="color:red">**!**</span> heterozygous data (phased vs. unphased!)

**Concatenation versus species tree inference (coalescent analyses)**

Concatenation:

NJ – Geneious

MP – TNT

ML – RAxML (…),  ~ FastTree

BI – MrBayes, BEAST …

# Coalescent-based Methods for Species Tree Inference

- **Summary statistic methods:** Start with estimated gene trees

  - ► Using estimated branch lengths:

    - ★ STEM (Kubatko et al. 2009)

    - ★ STAR, STEAC (Liu et al. 2009)

  - ► Using topology information only:

    - ★ Minimize Deep Coalescences (PhyloNet; Than & Nakhleh 2009)

    - ★ MP-EST (Liu et al. 2010)

    - ★ ST-ABC (Fan and Kubatko 2011)

    - ★ STELLS (Wu 2011)

- **Methods that utilize the full data:** Input is aligned sequences

  - ► BEST (Liu and Pearl 2007)

  - ► *BEAST (Heled and Drummond 2010)

  - ► New method based on algebraic statistics (Chifman and Kubatko 2013)

- Comparison of approaches:

  - Summary statistic methods

    - ★ Advantage: Quick

    - ★ Disadvantage: Ignore information in data

    - ★ Most current implementations do not easily allow for assessment of uncertainty

  - Full data methods

    - ★ Advantage: Fully model-based framework

    - ★ Disadvantage: Computationally intensive, sometimes prohibitively so

    - ★ Both BEST and *BEAST utilize a Bayesian framework and involve MCMC